

18/09

SOURCES OF APPROXIMATION

- 1) MEASURE ERRORS → measure instrument
- 2) ARITHMETIC ERRORS → propagation of rounding errors operation by operation
- 3) TRUNCATION ERRORS → truncation of infinite
- 4) INHERENT ERRORS → finite representation of data

(ABSOLUTE VS RELATIVE) ERROR

$$\left. \begin{array}{l} \text{ABSOLUTE : } E_x = \tilde{x} - x \\ \text{RELATIVE : } R_x = \frac{\tilde{x} - x}{x} \quad (x \neq 0) \end{array} \right\} \text{where } \tilde{x} \approx x \longrightarrow \tilde{x} = x(1 + r_x)$$

ACCURACY VS PRECISION , (TRUNCATION VS ROUNDING) ERROR NO!!!

REPRESENTATION OF REAL-NUMBER IN BASE β

HP: $\beta > 1$, $x \neq 0$, $0 \leq d_i \leq \beta - 1$, $d_1 \neq 0$ and not all $d_i = \beta - 1$

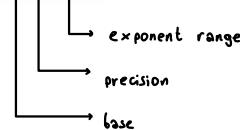
if not, you may have infinite representations of the same number by adding zeros

$0.d_1d_2d_3\dots \leftarrow \frac{1}{\beta} \leq m < 1 \leftarrow \text{mantissa}$

TH: $x = \text{sign}(x)(d_1\beta^{-1} + d_2\beta^{-2} + \dots)\beta^p = \text{sign}(x)m\beta^p \longrightarrow \text{UNIQUE REPRESENTATION:}$

FLOATING-POINT NUMBER

$\mathcal{F}(\beta, t, L, U)$ system of floating points



exponent of characteristics
 \uparrow
 $(d_1 \dots d_t) = m \longrightarrow \text{mantissa}$

The Floating point $x = \pm(d_1\beta^{-1} + \dots + d_t\beta^{-t})\beta^p$, $L \leq p \leq U$, $0 \leq d_i \leq \beta - 1$, $i = 1, \dots, t$

It is normalized when $d_1 \neq 0 \longrightarrow \frac{1}{\beta} \leq d_1 < 1 \longrightarrow ???$

- REASONS
- Representation of each number unique.
 - No digits wasted on leading zeros.
 - Leading bit need not be stored (in binary system).

[es]: $x = 12.\bar{3}$

$$\mathcal{F}(10, 5, -3, 3) = \pm(0.12333) \cdot 10^2$$

TRUNCATED $\leftarrow 3 < 5$

$x = 12.\bar{8}$

$$\mathcal{F}(10, 5, -3, 3) = \pm(0.12889) \cdot 10^2$$

$8 > 5 \rightarrow \text{ROUNDED}$

PROPERTIES OF $F(\beta, t, L, U)$

- 1) System is finite and discrete
- 2) TOT_NUM PTS = $2(\beta - 1)\beta^{t-1}(U - L + 1) + 1$
- 3) SMALLEST_NUM_POSITIVE = UFL = β^{L-1} → $F(\beta, t, L, U) \ll \rightarrow m \ll \text{ and } t \ll \frac{1}{\beta}$
 $\hookrightarrow UFL = \pm m \cdot \beta^t = \pm \beta^{-1} \cdot \beta^L = \beta^{L-1}$
- 4) LARGEST_NUM = OFL = $\beta^U(1 - \beta^{-t})$ → $F(\beta, t, L, U) \gg \rightarrow m \gg \text{ and } t \gg U$
 $\hookrightarrow OFL = \pm m \cdot \beta^t = \dots$
- 5) Numbers are equally spaced only between successive powers of β
- 6) Not all real numbers are exactly representable → only MACHINE NUMBERS are
 \hookrightarrow elements of $F(\beta, t, L, U)$

ROUNDING RULES

- 1) CHOP or TRUNCATION
- 2) ROUND TO NEAREST

STANDARD IEEE

- 1) SINGLE PRECISION → 4 byte = 32 bit = 1(sign) + 8(m) + 23(t) $F(2, 24, -128, 127)$
 - 2) DOUBLE PRECISION → 8 byte = 64 bit = 1(sign) + 11(m) + 52(t) $F(2, 53, -1024, 1023)$
- Ps: in base 2, if it is normalized ($d_1 \neq 0$) d_1 must be 1, so it's not stored to have more space → so mantissas are 23 (not 24) and 52 (not 53)

MACHINE PRECISION (ϵ_{mach})

If we represent numbers by:

- 1) CHOPPING → $\epsilon_{mach} = \beta^{1-t}$
- 2) ROUNDING TO NEAREST → $\epsilon_{mach} = \frac{1}{2}\beta^{1-t}$

DEF: ϵ_{mach} is smallest value of ϵ such that $f_l(1 + \epsilon) > 1$ and we could

also write it as
$$\left| \frac{f_l(x) - x}{x} \right| \leq \epsilon_{mach}$$

 $\hookrightarrow |R_x|$

EXCEPTIONAL VALUES

In IEEE we have 'Inf' and 'NaN'
 ↪ $1/0$ ↪ $0/0$; $0 \cdot Inf$; Inf/Inf

FLOATING POINT ARITHMETIC

$x \odot y = fl(x \cdot y)$ where $\odot: F \times F \rightarrow F$
which error is $\left| \frac{(x \odot y) - (x \cdot y)}{x \cdot y} \right| < \text{machine precision of the finite set of numbers we are using}$

STEPS TO PERFORM $x \odot y$

where $\{x \in F, y \in F, \odot \in \{+, -, \times, /\}\}$

- 1) $x \odot y$ in a register, which is a portion of memory which is double or larger than the portion of memory supposed to store x and y . Because when you perform operations you can have extra digits so you need more space to store them.
- 2) $fl(x \odot y) \in F$ we must compute the float of the result and this is stored in the finite set, getting truncated or rounded \rightarrow error

CANCELLATION

It's the subtraction between 2 t-digits numbers with same sign and similar magnitude.

It returns a result with less than t-digits \rightarrow EXACTLY REPRESENTABLE BUT WITH SERIOUS LOSS OF INFORMATION

es $x = 1.92403 \cdot 10^2, y = 1.92295 \cdot 10^2$

$$fl(x) = 0.192403 \cdot 10^3, fl(y) = 0.192295 \cdot 10^3$$

$$\text{Then } z = fl(x-y) = fl([1.92403 - 1.92295] \cdot 10^2) = fl(0.00128 \cdot 10^2) = 0.128 \\ fl(z) = 0.128000$$

19/09

LINEAR SYSTEMS

It's written as $Ax = b$ also written as

$$\begin{bmatrix} m \times n \\ A_{11} \dots A_{1n} \\ \vdots \dots \vdots \\ A_{m1} \dots A_{mn} \end{bmatrix} \cdot \begin{bmatrix} n \times 1 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} m \times 1 \\ b_1 \\ \vdots \\ b_m \end{bmatrix}$$

$\left\{ \begin{array}{l} A_{11}x_1 + A_{12}x_2 + \dots + A_{1n}x_n = b_1 \\ A_{m1}x_1 + A_{m2}x_2 + \dots + A_{mn}x_n = b_m \end{array} \right.$

We'll see:

- 1) EXISTENCE / UNIQUENESS OF SOLUTION
- 2) NUMERICAL METHODS FOR SOLUTION
- 3) CONDITIONING OF THE PROBLEM

SQUARE LINEAR SYSTEMS

 $\det(A) \neq 0$ and A invertible

RANK IS MAXIMUM

HP: $m=n$, (A not singular | $\text{rank}(A)=n$ | $Ax=0$ only for $x=0$)TH: $b = (b_1, \dots, b_n)$ exists and is unique

RESULT COMPUTATION

$$\begin{aligned} 1) \quad & Ax = b \rightarrow A^{-1}A x = A^{-1}b \rightarrow x = A^{-1}b && \text{EXPENSIVE IF } n \gg \\ & \text{COFACTOR MATRIX} \end{aligned}$$

$$2) \quad x_i = \frac{\det(A_i)}{\det(A)}$$

ALSO

• So we'll consider 2 methods.

- 1) DIRECT → FINITE STEPS × SOLUTION, MORE PRECISE, EXPENSIVE COST
- 2) ITERATIVE → INFINITE STEPS × SOLUTION, LESS PRECISE, LESS EXPENSIVE

• In order to perform error analysis:

- 1) ROUNDING/ARITHMETIC ERRORS × NUM_ALGORITHM_STEPS
(IT'S CALLED STABLE IF WE COULD SAY THAT TOTAL_ERROR < ε)
- 2) HIDDEN ERRORS ✗ NUM_ALGORITHM_STEPS
(IT'S CALLED ILL-POSED PROBLEM)
(IF ERROR >> ε)

DIRECT METHODS

From $Ax = b$ let $A = LU$ so that we have
 \downarrow
 $LUx = b \quad O\left(\frac{n^3}{3}\right)$

$$\begin{cases} Ly = b & \text{FORWARD SUBSTITUTION ALGORITHM} \\ Ux = y & \text{BACKWARD SUBSTITUTION ALGORITHM} \end{cases}$$

THE TOTAL COST IS
 $O\left(\frac{n^3}{3}\right) + 2 \cdot O(n^2)$
 DOMINANT

Where L is lower triangular matrix \rightarrow FORWARD RESOLUTION
 Where U is upper triangular matrix \rightarrow BACKWARD RESOLUTION

When we have pivoting algorithm $\rightarrow PA = LU \rightarrow PAx = Pb \rightarrow \begin{cases} Ly = Pb \\ Ux = y \end{cases}$

So we have 3 steps:

- 1) $PA = LU$
- 2) $Ly = Pb$
- 3) $Ux = y$

ITERATIVE METHODS

The idea is to build a vector using convergence $X^* = \lim_{K \rightarrow \infty} X_K$

$$\text{EXACT SOLUTION} \quad X^* = \lim_{K \rightarrow \infty} X_K \quad g(X_{K-1}) \text{ IN GENERAL}$$

We have different examples depending on their $g(X_{K-1})$:

1) STATIONARY ITERATIVE METHODS

$$X_{K+1} = B X_K + F \quad \text{NUMBER OF PERFORMED ITERATIONS} \quad \rightarrow \text{THE SMALLER IS } P, \text{ THE QUICKER IS THE ALGORITHM}$$

ITERATION MATRIX VECTOR FROM b

2) GRADIENT LIKE METHODS

SKIPPATA

We cannot iterate forever so we must know when to stop, in order to do so there are STOPPING CRITERIA:

- 1) $\|r_K\| \leq \epsilon$
 - 2) $\frac{\|r_K\|}{\|b\|} \leq \epsilon$
 - 3) $\|X_{K+1} - X_K\| \leq \tau$
- $r_K = b - Ax_K$

$$\left. \begin{array}{l} \text{TOLERANCE } \approx [10^{-5}, 10^{-10}] \\ \|r_K\| \leq \epsilon \end{array} \right\}$$

ARRAY NORMS

$$\|\cdot\| : \mathbb{R}^n \rightarrow \mathbb{R}$$

- 1) $\|x\| \geq 0$ and $\|x\| = 0 \Leftrightarrow x = 0$
- 2) $\|x+y\| \leq \|x\| + \|y\|$
- 3) $\|\lambda x\| = |\lambda| \cdot \|x\|$

EUCLIDEAN NORM :

$$\|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$$

MANHATTAN NORM :

$$\|x\|_1 = \sum_{i=1}^n |x_i|$$

INFINITE NORM :

$$\|x\|_\infty = \max_{1 \leq i \leq n} \{|x_i|\}$$

MATRIX NORMS

$$\|\cdot\| : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$$

- 1) $\|A\| \geq 0$ and $\|A\| = 0 \Leftrightarrow A = 0$
- 2) $\|A+B\| \leq \|A\| + \|B\|$
- 3) $\|c \cdot A\| = |c| \cdot \|A\|$

FORBE NIUS NORM :

$$\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n a_{ij}^2}$$

2-NORM:

$$\|A\|_2 = \sqrt{\rho(A^T A)}$$

THE LARGEST EIGEN VALUES IN ABSOLUTE VALUE

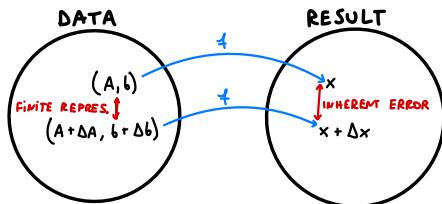
AUTOVALORE

NON È MOLTIPLICAZIONE MA INDICA DI A CHI SI RIFERISCE

1-NORM:

$$\|A\|_1 = \max_{1 \leq i \leq n} \left\{ \sum_{j=1}^m |a_{ij}| \right\}$$

INHERENT ERROR



When the error in the result is much greater than the error in the data we have a problem, otherwise it's normal.

So the problem is when $\|\Delta x\| \gg \|\Delta A\| + \|\Delta b\|$ PROBLEM IS ILL-CONDITIONED

RESULT :
$$\left\| \frac{\Delta x}{x} \right\| \leq \|A^{-1}\| \cdot \|A\| \cdot \frac{\|\Delta b\|}{\|b\|} = K(A) \frac{\|\Delta b\|}{\|b\|}$$

RELATIVE MEASURE OF THE ERROR

RELATIVE MEASURE OF THE DATA

CONDITION NUMBER OF A MATRIX

$K(A)$ small \rightarrow WELL-CONDITIONED

$K(A)$ large \rightarrow ILL-CONDITIONED

21/09

DOT PRODUCT

$$f : \mathbb{R}^n \times \mathbb{R}^n \longrightarrow \mathbb{R}$$

$$x^\top y = \sum_{i=1}^n x_i y_i$$

SYMMETRIC, POSITIVE DEFINITE MATRICES

It's a matrix $A \in \mathbb{R}^{n \times n}$ where :

- 1) $A = A^\top \longrightarrow$ SYMMETRIC
- 2) $\forall x \in V - \{\emptyset\} : x^\top A x > \emptyset \longrightarrow$ POSITIVE SEMI DEFINITIVE
- 3) $\forall x \in V - \{\emptyset\} : x^\top A x > 0 \longrightarrow$ DEFINITE POSITIVE

And it follows those properties:

- 1) $\text{Ker}(A) = \{\emptyset\}$ where $\text{Ker}(A) := \{x \in \mathbb{R}^n : Ax = \emptyset\}$
Because $x^\top A x > \emptyset \forall x \neq \emptyset \longrightarrow Ax \neq \emptyset \text{ if } x \neq \emptyset$
- 2) $a_{ii} > \emptyset$
Because $a_{ii} = e_i^\top A e_i > \emptyset$ λ -th vector of standard basis in \mathbb{R}^n

- 3) If the eigenvalues of $A \in \mathbb{R}^{n \times n}$ symmetric are all positive, then A is positive definite

ANGLES AND ORTOGONALITY

The scalar product can provide us information about the angles between 2 vectors :

$$\cos \omega = \frac{\langle x, y \rangle}{\|x\| \|y\|}.$$

If $\omega < \pi$ then x, y have the same orientation and vice-versa.

ORTOGONALITY: x, y are orthogonal $\longleftrightarrow \langle x, y \rangle = \emptyset$

ORTONORMALITY: x, y are orthogonal and $\|x\| = 1 = \|y\|$

ORTOGONAL MATRIX

It's a matrix $A \in \mathbb{R}^{n \times n}$ where its column are orthonormal $\longrightarrow AA^\top = I = A^\top A \longrightarrow A^{-1} = A^\top$

- 1) $\|Ax\|^2 = \|x\|^2 \longrightarrow$ the length of a vector x is not changed when transforming it by orthogonal matrix.
- 2) $\cos(\omega) = \frac{x^\top y}{\|x\| \|y\|} \longrightarrow$ the angle between Ax and Ay is the angle between x and y

ORTHONORMAL BASIS

$$c_1v_1 + \dots + c_nv_n = 0 \iff c_1 = \dots = c_n = 0$$

Giving a vector space "V" of dimension n, $\{v_1, \dots, v_n\}$ linear independent is a basis of V
DEF: $\{b_1, \dots, b_n\}$ is a basis of V if $\begin{cases} \langle b_i, b_j \rangle = 0 & \text{for } i \neq j \\ \langle b_i, b_i \rangle = 1 \end{cases}$

ORTHONORMAL COMPLEMENT

If V vector space and M vector space : $M \subseteq V$ then M is a subspace.
 Then there is M^\perp which is a V subspace which contains all vector in V that are orthogonal to every vector in M and $M \cap M^\perp = \{0\}$

ORTHOGONAL PROJECTIONS

They are important in machine learning because of DATA REDUCTION.
 We have to deal with high-dimensional data which is hard to analyze.

WE NEED TO CHANGE THE DATA \rightarrow COMPRESSION

- We project the original high-dim data in low-dim space to learn more about dataset and extract relevant patterns

DEF: $U \subseteq V$, a linear mapping $\pi: V \rightarrow U$ it's an ORTHOGONAL PROJECTION if $\pi \circ \pi = \pi$

can be expressed by transformation matrices \rightarrow PROJECTION MATRICES $\rightarrow P_\pi^2 = P_\pi$

PROJECTION ONTO GENERAL SUBSPACES

Let $V = \mathbb{R}^n$, so we look at vectors $x \in \mathbb{R}^n \geq U$, where $\dim(U) = m \geq 1$

Assume (b_1, \dots, b_m) an ordered basis of U, so $\pi_U(x)$ is any projection of x onto U and so it should be necessarily be an element of U.

$$\pi_U(x) = \sum_{i=1}^m \lambda_i b_i \rightarrow \text{LINEAR COMBINATION OF } (b_1, \dots, b_m)$$

EIGENVALUES AND EIGENVECTORS

Let $A \in \mathbb{R}^{n \times n}$, $\lambda \in \mathbb{C}$ is EIGENVALUES of A and $x \in \mathbb{R}^n - \{0\}$ is the EIGENVECTOR if $Ax = \lambda x$

We have 4 equivalent statements

1) λ is EIGENVALUE of $A \in \mathbb{R}^{n \times n}$

2) $\exists x \in \mathbb{R}^n - \{0\} : Ax = \lambda x$

3) $(A - \lambda I_n)x = 0$ with $x \neq 0$

4) $\text{rk}(A - \lambda I_n) < n$

5) $\det(A - \lambda I_n) = 0$

THE EIGENVECTORS (x_1, \dots, x_n) OF $A \in \mathbb{R}^{n \times n}$ WITH n DIFFERENT $(\lambda_1, \dots, \lambda_n)$ ARE LINEARLY INDEPENDENT

• SYMMETRIC, POSITIVE DEFINITE MATRICES \rightarrow POSITIVE, REAL EIGENVALUES

• If x is an eigenvector of A associated with eigenvalue λ , also cx will be an eigenvector of A with the same eigenvalue

EIGENSPACE AND EIGENPECTRUM

EIGENSPACE: the subspace of \mathbb{R}^n spanned by the set of eigenvectors of A with eigenvalues λ

EIGENPECTRUM: the set of all eigenvalues of A

GEOMETRIC MULTIPLICITY EIGENVALUE

It's the number of linearly independent eigenvectors associated with λ_i .

DEFECTIVE MATRIX

A square matrix $A \in \mathbb{R}^{n \times n}$ is defective if it possesses fewer than n linearly independent eigenvectors.

SPECTRAL THEOREM

HP: $A \in \mathbb{R}^{n \times n}$ symmetric

TH: \exists an orthonormal basis of the corresponding vector space V consisting of eigenvectors of A and each eigenvalue is real.

IN OTHER WORDS: If A is symmetric, then the matrix, containing its eigenvectors, is orthogonal and the eigenvalues are real

SIMILAR MATRICES

A, D are similar if exists an invertible one P so that $D = P^{-1}AP$

- A, D so will have the same eigenvalues

DIAGONALIZABLE MATRIX

$A \in \mathbb{R}^{n \times n}$ is diagonalizable if it's similar to a diagonal matrix $\longrightarrow D = P^{-1}AP$

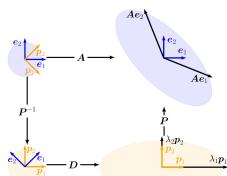
- Symmetric matrix can always be diagonalized

EIGENDECOMPOSITION

$A \in \mathbb{R}^{n \times n}$ can be factored into $A = PDP^{-1}$ where D is a diagonal matrix with only eigenvalues of A on the diagonal

eigenvectors of A form a basis of \mathbb{R}^n

- Geometrically we can see matrix as a linear transformation between 2 vector spaces



25/09

SINGULAR VALUE COMPOSITION (SVD)

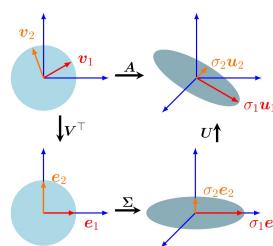
- More general theorem → applied to all matrices and it's always possible

HP: $A^{m \times n}$ and $\text{rank}(A) = r \in [0, \min\{m, n\}]$

TH:

$$\begin{matrix} & n \\ \begin{matrix} m \\ A \end{matrix} & = \begin{matrix} m \\ U \end{matrix} \begin{matrix} n \\ \Sigma \end{matrix} \begin{matrix} n \\ V^T \end{matrix} \end{matrix} \xrightarrow{\text{ORTHOGONAL MATRIX}}$$

ORTHOGONAL MATRIX



$\sum_{ii} = \sigma_i$, $\sum_{ij} = 0$ for $i \neq j$ → CAN'T SPEAK PROPERLY ABOUT DIAGONAL MATRIX ($m \times n$)

BY CONVENTION
 $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r \geq 0$

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \sigma_n \\ 0 & \dots & 0 \\ \vdots & & \vdots \\ 0 & \dots & 0 \end{bmatrix}$$

SINGULAR VALUES

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & 0 & 0 & \dots & 0 \\ 0 & \ddots & 0 & 0 & \dots & 0 \\ 0 & 0 & \sigma_m & 0 & \dots & 0 \\ 0 & 0 & 0 & \ddots & \dots & 0 \\ \vdots & & & & & \vdots \\ 0 & \dots & 0 & 0 & \dots & 0 \end{bmatrix}$$

 $m - n$

PROPERTIES OF SVD

- There's no relation between eigenvalues and singular values of A , but there is between the singular values and the eigenvalues of $A^T A$

HP: $\lambda_i \geq 0$ eigenvalues of $A^T A$

TH: $A V_i = \sigma_i u_i$

DiM:

$$(A \cdot B)^T = B^T A^T \xrightarrow{\text{I}}$$

$$A^T A = (U \Sigma V^T)^T (U \Sigma V^T) = V \Sigma^T U^T U \Sigma V^T = V \Sigma^T \Sigma V^T = V \Sigma^2 V^T$$

$$A A^T = (U \Sigma V^T) (U \Sigma V^T)^T = U \Sigma V^T V \Sigma^T U^T = U \Sigma \Sigma^T V^T = U \Sigma^2 V^T$$

Then we can obtain the SINGULAR VALUE EQUATION: $A = U \Sigma V^T \longrightarrow AV = U \Sigma V^T$

$$AV_i = \sigma_i u_i$$

EIGENVALUE DECOMPOSITION VS SINGULAR VALUE DECOMPOSITION

EIGENVALUE DECOMPOSITION VS SINGULAR VALUE DECOMPOSITION

1) $\exists P \in \mathbb{R}^{n \times n}$

2) Vectors in P are not necessarily orthogonal

NOT NECESSARILY A ROTATION

1) $\exists U, V \in \mathbb{R}^{m \times n}$

2) U and V are orthonormal

They are both compositions of three linear mappings → REPRESENTS A ROTATION

$$\text{EIGENVALUES OF } A \leftarrow \lambda_i = \sigma_i^{-2} \rightarrow \text{SINGULAR VALUES OF } A$$

RELATION BETWEEN SINGULAR VALUES AND $\|A\|_2$

$$\|A\|_2 = \sqrt{\max\{|\lambda_1|, \dots, |\lambda_n|\} (A^T A)}$$

\downarrow

$\sigma_1, \dots, \sigma_n$ are sorted

$$\max\{\sigma_1, \dots, \sigma_n\}(A) = \sigma_1(A)$$

$$\|A^{-1}\|_2 = \sqrt{\rho((A^{-1})^T A^{-1})} = \sqrt{\rho((AA^T)^{-1})} = \frac{1}{\sigma_n(A)}$$

IF $\lambda_1, \dots, \lambda_n$ eigenvalues of A

THEN $\frac{1}{\lambda_1}, \dots, \frac{1}{\lambda_n}$ eigenvalues of A^{-1}

$$K_2(A) = \frac{\sigma_1(A)}{\sigma_n(A)}$$

MATRIX APPROXIMATION

We want to use SVD in a cheaper way than its own computation

HP: $A = U \sum V^T \in \mathbb{R}^{m \times n}$

TH: $\hat{A}(k) := \sum_{i=1}^k \sigma_i u_i v_i^T = \sum_{i=1}^k \sigma_i A_i \quad \simeq A$

DIM: $A_i = u_i \cdot v_i^T \quad (u_i \in \mathbb{R}^m, v_i \in \mathbb{R}^n)$

From HP: $A = \sum_{i=1}^r \sigma_i u_i v_i^T = \sum_{i=1}^r \sigma_i A_i$

PS: $\hat{A}(k)$ provides $K \cdot (1+m+n)$ numbers

$\hat{A}(k) = \sum_{i=1}^k \sigma_i A_i \quad \simeq A$

$K < r$

28/09

SPECTRAL NORM OF A MATRIX

For $x \in \mathbb{R}^n - \{0\}$, the spectral norm of $A \in \mathbb{R}^{m \times n}$ is $\|A\|_2 := \max_x \frac{\|Ax\|_2}{\|x\|_2}$

DETERMINES HOW LONG ANY VECTOR x
CAN BECOME WHEN MULTIPLIED BY A

SPECTRAL NORM AS APPROXIMATION

HP: $A \in \mathbb{R}^{m \times n}$ $\text{rank}(A) = r$, $B \in \mathbb{R}^{m \times n}$ $\text{rank}(B) = k$

TH: $\forall k: k \leq r$, $\hat{A}(k) = \sum_{i=1}^k \sigma_i u_i v_i^T$ we have that $\hat{A}(k) = \underset{\text{THE MINIMUM OVER THE SET OF ALL MATRICES OF } m \times n \text{ WITH RANK } k}{\text{argmin}_{\text{rk}(B)=k}} \|A - B\|_2$

$$?? \quad 2) \|A - \hat{A}(k)\|_2 = \sigma_{k+1}$$

THE MINIMUM OVER THE SET OF ALL MATRICES OF $m \times n$ WITH RANK k
REPRESENTING THE MATRIX WHICH HAS
MINIMUM DISTANCE IN RESPECT TO A INTO
2-NORM

MACHINE LEARNING (LECTURE 4 - EVANGELISTA SITE)

- Machine Learning is the set of all techniques and algorithms able to extract knowledge from the data and use it to make accurate decisions

The development of a Machine Learning algorithm follows some steps:

1 - UNDERSTANDING

A task is an unknown function $f(x)$. We need to approximate $y = f(x) \approx f_0(x)$

- 1) is it learnable? \rightarrow if there's a real correlation between input x and output y
- 2) is it collectable? \rightarrow if we can collect data x efficiently and with scalability

2 - COLLECTION

Data requires time to be collected in good quality, usually it's necessary to clean up the dataset with some specific algorithm

We take data in .csv, we convert strings into numbers and then we collect all these numbers in a big matrix and we know how to manipulate them.

We have 2 way to encode str \rightarrow int

- 1) INTEGER ENCODING \rightarrow assign integer to value, summing doesn't make sense and creates undesired relations
- 2) ONE-HOT ENCODING \rightarrow assign vector of integers to value, no unwanted relations and summing makes sense

After that the result is $X = [x^1 \ x^2 \ \dots \ x^N] \in \mathbb{R}^{d \times N}$ (d -features, N -num of datapoints)

$$y^i = f(x^i) \quad Y = [y^1 \ y^2 \ \dots \ y^N]^T \in \mathbb{R}^N$$

3 - DESIGN

- 1) SUPERVISED LEARNING → dataset composed by inputs x and outputs y → CLUSTERING
- 2) UNSUPERVISED LEARNING → dataset composed by only inputs x

4 - TRAINING

We give data to let our model learn the patterns

5 - TUNING

It's the task of finding the good compromise over the flexibility of the model depending on the task.

OVERFITTING UNDERFITTING

6 - TESTING

$$N_{\text{DATA_TOT}} = N_{\text{DATA_TEST}} + N_{\text{DATA_TRAIN}} \rightarrow \text{SPLITTING}$$

In order to maximize the accuracy of the prediction of your model :

$$\text{Acc}(f_\theta) = \frac{1}{N_{\text{test}}} \sum_{i=1}^{N_{\text{test}}} \ell(f_\theta(x_i), k_{x_i})$$

Diagram illustrating the components of the accuracy function:

- MODEL FUNCTION**: Points to f_θ .
- ACCURACY FUNCTION**: Points to ℓ .
- MODEL OUTPUT**: Points to k_{x_i} .
- REAL OUTPUT**: Points to x_i .

28/09

LINEAR LEAST SQUARE PROBLEM (LSQ)

If we have $Ax = b$ we can't compute $Ax - b = 0$ but we can find \tilde{x} : $A\tilde{x} - b \approx 0$, so we have to find the smallest r possible $\rightarrow \min \{ \|r\|_2\}$

$$\text{So we have } \tilde{x} = \min_{x \in \mathbb{R}^n} \|r\|_2 = \min_{x \in \mathbb{R}^n} \|r\|_2^2 \quad \text{WE ELIMINATE THE SQUARE ROOT}$$

I DON'T FORGET THAT YOU ARE
LOOKING FOR ARG()

HP: $A \in \mathbb{R}^{m \times n}$, $m > n$, $r = \text{rank}(A)$

TH: 1) $r = n \rightarrow$ unique solution

2) $r < n \rightarrow$ infinite solutions

Dim: 1) $r = n$
 $\begin{array}{c} n \times m \text{ symmetric positive} \\ \text{definite matrix} \\ \uparrow \quad \downarrow \\ A^T A x = A^T b \quad \text{if rank}(A) = n \end{array}$
 $Ax = b \rightarrow A^T A x = A^T b \rightarrow \text{UNIQUE SOLUTION} \rightarrow$ we can use direct methods to compute it

2) $r < n$

$$Ax^* = b \rightarrow S = \{x \in \mathbb{R}^n : x \text{ solutions of } \min_{x \in S} \|Ax - b\|_2^2\}$$

$$\xrightarrow{\quad} x^* = \min_{x \in S} \|x\|_2$$

WE USE SVD
DECOMPOSITION

$$x^* = \sum_{i=1}^r \frac{b_i}{\sigma_i} \begin{pmatrix} M_i^T \\ V_i \end{pmatrix}$$

APPLICATION OF LSQ

$$f(x) = c_0 + c_1 x + c_2 x^2 + \dots + c_n x^n$$

$$\Theta = \{c_0, c_1, \dots, c_n\}$$

$$r = \{r_0, r_1, \dots, r_n\} = \{y_0 - f(x_0), \dots, y_n - f(x_n)\} = \{y_0 - (c_0 + c_1 x_0 + \dots + c_n x_0^n), \dots, y_n - (c_0 + c_1 x_n + \dots + c_n x_n^n)\}$$

So we have $\vec{r} = \vec{y} - A\vec{x}$

PARTIAL DIFFERENTIATION

Given $f(x) = (f_1(x), \dots, f_n(x))$ then:

$$1) \text{ GRADIENT} \longrightarrow \nabla f(x) = \left(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n} \right) \xrightarrow{\text{CHAIN - RULE}} \frac{\partial(g \circ f)}{\partial x} = \frac{\partial g}{\partial f} \cdot \frac{\partial f}{\partial x}$$

2) HESSIAN MATRIX

$$\mathbf{H}_f = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$$

3) JACOBIAN MATRIX

$$\mathbf{J}_f = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

GRADIENT OF A
MULTIVARIABLE
FUNCTION

LEAST-SQUARES FUNCTION

We consider a linear model $y = \Phi(\theta) \rightarrow \mathbb{R}^D$ PARAMETER
 OBSERVATIONS \mathbb{R}^N INPUT FEATURES $\mathbb{R}^{N \times D}$

Then we consider $L(\theta) = \|y - \Phi(\theta)\|_2^2$ as LEAST-SQUARES FUNCTION

Now we look for $\frac{\partial L}{\partial \theta} \in \mathbb{R}^{1 \times D}$:

$$\frac{\partial L}{\partial \theta} = \underbrace{\frac{\partial L}{\partial (y - \Phi(\theta))}}_1 \cdot \underbrace{\frac{\partial (y - \Phi(\theta))}{\partial \theta}}_2 = 2(y - \Phi(\theta))^T(-\Phi) = -2(Y^T - \Phi^T(\theta^T)\Phi)$$

$$1) \frac{\partial \|y - \Phi(\theta)\|_2^2}{\partial (y - \Phi(\theta))} = \frac{\partial \|e\|_2^2}{\partial e} = \left[\frac{\partial}{\partial e_1} \left(\sqrt{\sum_{k=1}^n e_k^2} \right)^2, \dots, \frac{\partial}{\partial e_n} \left(\sqrt{\sum_{k=1}^n e_k^2} \right)^2 \right] = \left[\frac{\partial}{\partial e_1} \sum_{k=1}^n e_k^2, \dots, \frac{\partial}{\partial e_n} \sum_{k=1}^n e_k^2 \right] =$$

$$= \left[\sum_{k=1}^n \frac{\partial e_k^2}{\partial e_1}, \dots, \sum_{k=1}^n \frac{\partial e_k^2}{\partial e_n} \right] = [2e_1, \dots, 2e_n] = 2e^T = 2(Y - \Phi(\theta))^T$$

$$2) \frac{\partial (y - \Phi(\theta))}{\partial \theta} = -\Phi$$

BACK PROPAGATION AND AUTOMATIC DIFFERENTIATION

We can't use chain-rule explicitly for computing the gradient, because otherwise we would introduce unnecessary overhead.

GRADIENT IN A DEEP NETWORK

The chain-rule is massively used in deep-learning, where we have much more function composed.

So we use BACKPROPAGATION

$$\frac{\partial L}{\partial \theta_{k-1}} = \frac{\partial L}{\partial f_k} \cdot \frac{\partial f_k}{\partial \theta_{k-1}} = \frac{\partial L}{\partial f_k} \cdot \frac{\partial f_k}{\partial \theta_{k-1}} \cdot \dots \cdot \frac{\partial f_{i+2}}{\partial f_{i+1}} \cdot \frac{\partial f_{i+1}}{\partial \theta_i}$$

SOLVE THE LAST ONE AND REUSE IT TO SOLVE
THE PREVIOUS ONE UNTIL WE SOLVED THE 1^o

AUTOMATIC DIFFERENTIATION

Backpropagation is a generalization of AUTOMATIC DIFFERENTIATION, which is a set of techniques to evaluate the gradient of function.

We use intermediate variables to handle the backpropagation of the chain-rule and we use the associativity of a matrix multiplication to handle the 2 directions of backpropagation

$$\frac{dy}{dx} = \frac{dy}{db} \cdot \left(\frac{db}{da} \cdot \frac{da}{dx} \right) \quad \text{FORWARD MODE} \rightarrow \text{gradients are expanded}$$

$$\frac{dy}{dx} = \left(\frac{dy}{db} \cdot \frac{db}{da} \right) \cdot \frac{da}{dx} \quad \text{REVERSE MODE} \rightarrow \text{isolate a gradient from the one before it}$$

Let (x_1, \dots, x_d) input values

(x_{d+1}, \dots, x_D) intermediate variables

x_D output variable

OPTIMIZATION

Procedure to find max/min of function. It's important in machine learning because when we want to find the parameters during the training part, we have to maximize/ minimize a suitable function in order to find the best value.

CONDITION FOR \exists global/local MINIMUM

$f: \mathbb{R}^n \rightarrow \mathbb{R}$ continuous and differentiable in \mathbb{R}^n

WE DON'T CARE MAXIMUM

if all partial derivatives exist and are continuous

WITHOUT THE EQUAL IN THE DEFINITION WE TALK ABOUT THE 'strict' VERSION

- LOCAL MINIMUM $x^* \in \mathbb{R}^n$ is a local minimum $\iff f(x^*) \leq f(x) \quad \forall x: \|x - x^*\| \leq \epsilon$
- GLOBAL MINIMUM $x^* \in \mathbb{R}^n$ is a global minimum $\iff f(x^*) < f(x) \quad \forall x \in \mathbb{R}^n$
- STATIONARY POINT $x^* \in \mathbb{R}^n$ is a stationary point $\iff \nabla f(x^*) = \emptyset$

FIRST ORDER CONDITION

If x^* is a stationary point of local minimum for f , then $\nabla f(x^*) = \emptyset$ and x^* is a stationary point for f

$$x^* \text{ local minimum} \Rightarrow \nabla f(x^*) = \emptyset$$

SECOND ORDER CONDITION

f is twice differentiable $\rightarrow f$ is differentiable and also ∇f

If $\nabla f(x^*) = \emptyset$ and $\nabla^2 f(x^*)$ is positive definite, then x^* is a strict local minimum for f

$$x^* \text{ is a strict local minimum} \Leftrightarrow \begin{aligned} \nabla f(x^*) &= \emptyset \\ \nabla^2 f(x^*) &\text{ pos. def.} \end{aligned}$$

ALGORITHMS

They compute only local minimum

1) ITERATIVE METHODS

$$\vec{x}, \vec{x}_1, \vec{x}_2, \dots, \vec{x}_K \xrightarrow[k \rightarrow \infty]{\text{CONVERGES COMPONENT BY COMPONENT}} x^*$$

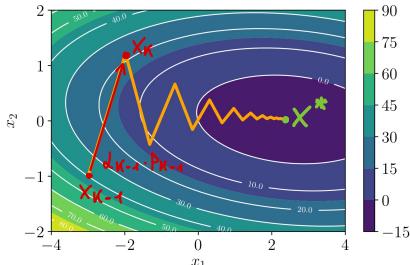
so $x_K = g(x_{K-1}, x_{K-2}) \rightarrow x_K \approx x^*$

DECENT METHOD \longleftrightarrow $g(x_K) < g(x_{K-1})$

$G?$

p_k : SEARCH DIRECTION

d_k : STEP LENGTH



CHOICE OF p

$\forall x \dots$

If $p \in \mathbb{R}^n$, $p \neq 0$ and $p^\top \nabla f(x) < 0$, then p is a decent direction for f in x

• DESCENT DIRECTION

$p \in \mathbb{R}^n$ is called descent direction for f in $x \longleftrightarrow \exists \lambda^* > 0 : f(x + \lambda p) < f(x) \quad \forall \lambda \in [0, \lambda^*]$

2) GRADIENT METHOD

$$x_{K+1} = x_K - d_K \nabla f(x_K)$$

CHOICE OF d_K

d_K is constant? \rightarrow NO \rightarrow the convergence is not guaranteed

BACKTRACKING PROCEDURE

$$d_K = 1$$

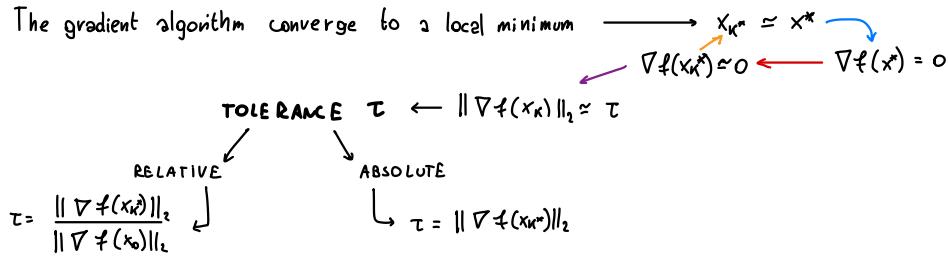
IF $f(x_K - d_K \nabla f(x_K)) > f(x_K) + \underbrace{\text{now se lo ricorda}}_{\text{INFORMATIC NOTATION}}$

$$\xrightarrow{?} f(w_t) - f(w_t + \lambda g_t) \geq \mu \lambda \|g_t^\top [\nabla f(w_t)]\|$$

Then λ is not small enough, so $d_K' = d_K / C$

Then algorithm converges to a local minimum

HOW TO CHOOSE K* WHERE TO STOP THE ALGORITHM



PSEUDOCODE OF GRADIENT METHOD

INPUT: f, x_0

$k = 0$

while [condition] :

```

p_k = - \nabla f(x_k)
a_k backpropagating procedure
x_(k+1) = x_k + a_k * p_k
k += 1
return x_k

```

BACKTRACKING PROCEDURE

OS/10

INITIALIZATION

In a sequence $x_{k+1} = x_k - \alpha_k \nabla f(x_k)$ we have $x_1 = x_0 - \alpha_0 \nabla f(x_0)$ at the beginning.

So x_0 is given in input and it's the seed to determine the whole sequence.

But how the start point can be chosen? $\rightarrow \text{random}(-1, 1)$ \rightarrow may unfit the problem for different reasons
 ↓
 DIFFERENT x_0 LEAD TO DIFFERENT x^* \leftarrow CHANGE THE INITIALIZATION POINT

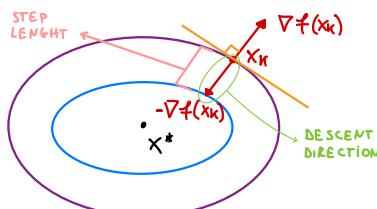
ALGORITHM PROBLEMS

1) FLAT REGIONS AND LOCAL OPTIMA



It can lead to slowing down the algorithm or at worse it could even stop it.
 But we can't know the shape before applying, so it occurs during execution.

2) DIFFERENTIAL CURVATURE



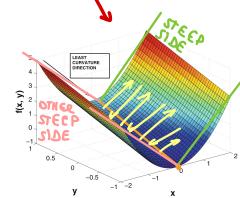
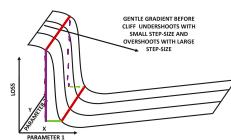
We draw the tangent to x_k in the level curves to a function.
 The gradient would be perpendicular in that point and going out of the curves. Speed of descent depends on gradient magnitude.
 THEY ARE THE STEEPEST MOVEMENT BUT ONLY IN THAT POINT

LEARNING RATE: $\|x_{k+1} - x_k\| = \alpha_k \|\nabla f(x_k)\| \rightarrow$ step length

If step lenght is too high (if α_k is too large) the method could diverge

We may have:

CLIFFS and VALLEYS



The gradient-descent will bounce along the steep sides, without making much progress.

We also have the VANISHING AND EXPLODING GRADIENT PROBLEM

It occurs when we have large differences in the magnitudes of the partial derivatives with respect to parameters of different layers.

A solution to this problem is to standardize features before applying gradient descent.

3) NON-DIFFERENTIABLE OBJECTIVE FUNCTIONS

Usually you know if a function it's differentiable or not, but if the function is implemented by a neural network we don't. NO DIFFERENTIABLE \rightarrow NO GRADIENT

If only a few points are not differentiable we can slightly modify the computation of the gradient and going on.

If there are more non-differentiable points we have serious problems in the computation of the gradient

CONVEX FUNCTIONS

CONVEX SET



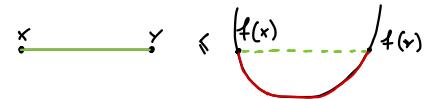
NON CONVEX SET



$$f: \Omega \subset \mathbb{R}^n \rightarrow \mathbb{R}$$

If Ω is a convex set and $\forall x, y \in \Omega$, $\forall \theta: 0 < \theta < 1$

$$\text{Then } f(\theta x + (1-\theta)y) \leq \theta f(x) + (1-\theta)f(y)$$



QUADRATIC FUNCTION

SIMMETRIC AND DEFINITE POSITIVE

$$f(\vec{x}) = \frac{1}{2} \vec{x}^T A \vec{x} + \vec{x}^T \vec{b} + c$$

$$\|\vec{A}\vec{x} - \vec{b}\|_2^2 = (\vec{A}\vec{x} - \vec{b})^T (\vec{A}\vec{x} - \vec{b}) = \vec{x}^T \vec{A}^T \vec{A} \vec{x} - \vec{b}^T \vec{A} \vec{x} - \vec{x}^T \vec{A} \vec{b} - \vec{b}^T \vec{b}$$

\Rightarrow so $\|\vec{A}\vec{x} - \vec{b}\|_2^2$ is quadratic function

PROPERTIES OF (STRICTLY) CONVEX FUNCTIONS

1) If f is convex

Then any point of local minimum is also a global minimum \rightarrow Δ there may be more global minimum

2) If f is strictly convex

Then it exists a unique point of global minimum

3) If f is convex and differentiable

Then any stationary point is a global minimum for f .

GRADIENT DESCENT WITH MOMENTUM

$$x_{i+2} = x_i - \alpha \nabla f(x_i) + \beta \Delta x_i$$
$$\Delta x_i = x_i - x_{i-2} \quad \text{--- SO WE MEMORIZE THE PREVIOUS 2 ITERATIONS}$$

STOCHASTIC GRADIENT DESCENT (SGD)

IN NEURAL NETWORKS

When we work with $N \gg$ we introduce much effort on gradient computation so we introduce an approximation in each step of gradient descent, called SGD.

STOCHASTIC \neq DETERMINISTIC

LOSS FUNCTION NUM. DATA ELEMENTS

$$L(\theta) = \sum_{n=1}^N L_n(\theta) \quad \rightarrow \quad \nabla L(\theta) = \sum_{n=1}^N \nabla L_n(\theta)$$

In practice the approximation is $\nabla L(\theta) \approx \sum_{n \in S} \nabla L_n(\theta) \quad S \subseteq \{1, \dots, N\}$ MINIBATCH

- When the learning rate decreases at an appropriate rate, stochastic gradient descent converges almost surely to local minimum.

EPOCH

Iterations necessary to use all the data. \rightarrow 1 ITERATION ALL OVER THE DATA

ES:

EPOCH	$i=1$	$\nabla L(\theta) \approx \nabla L_{n_1}(\theta) \quad n_1 \in \{1, \dots, N\}$
	$i=2$	$\nabla L(\theta) \approx \nabla L_{n_2}(\theta) \quad n_2 \in \{1, \dots, N\} - \{n_1\}$
	\vdots	\vdots
	$i=p$	$\nabla L(\theta) \approx \nabla L_{n_p}(\theta) \quad n_p \in \{1, \dots, N\} - \{n_1, \dots, n_{p-1}\}$

\downarrow 1 ELEMENT \rightarrow EPOCH FINISHED

09/10

STATE SPACE (Ω)

It's the set of all the possible results of an experiment

ES: two coins $\rightarrow \Omega = \{TT, HT, TH, HH\}$

EVENT (A)

It's a subset of Ω ($A \subset \Omega$) $\rightarrow \omega$ is particular event ($\omega \subset \Omega$)

PROBABILITY (P)

set of all possible result

It's a function $P(A) : S(A) \rightarrow [0, 1] \in \mathbb{R} \quad \forall A$

$$P(A) = \frac{\#(A)}{\#(\Omega)} \quad \xrightarrow{\text{CARDINALITY}} \text{EQUALLY LIKELY MODEL}$$

CONDITIONAL PROBABILITY

If B, A are dependent, $P(B|A) = \frac{P(A \cap B)}{P(A)}$ where $P(A) \neq 0$

$\hookrightarrow P(A \cap B) \neq P(A) \cdot P(B)$

\hookrightarrow EXTENSIBLE TO A_1, \dots, A_k EVENTS

RANDOM VARIABLES (X)

We need to map the events into mathematical objects.

$$X : \Omega \rightarrow \mathbb{R}$$

\uparrow \downarrow

$$\omega \quad X(\omega) = x$$

TARGET SPACE: $T = \{x : x = X(\omega), \omega \in \Omega\}$ $\xrightarrow{\text{SET OF ALL VALUES OF } X}$

Random variable can be DISCRETE or CONTINUOUS

DISCRETE RV

T is constituted by a finite or $\xrightarrow{\infty}$ number with precedent and successive number

of elements

CONTINUOUS RV

T is an interval or the union of intervals

DISCRETE RV CASE

PROBABILITY MASS FUNCTION (PMF) \longrightarrow distribution

$$f_x : T_x \rightarrow [0, 1]$$

$$\forall x \in T_x, f_x(x) = P(X=x)$$

PROPERTIES

$$1) f_x(x) \geq 0 \quad \forall x \in T_x$$

$$2) \sum_{x \in T_x} f_x(x) = 1$$

$$3) A \subset \Omega, P(X=x \in A) = \sum_{x \in A} f_x(x)$$

EXAMPLES

1) UNIFORM DISTRIBUTION : If $\#(T_x) = N$, $f_x(x) = 1/N$

2) POISSON DISTRIBUTION : If λ is the mean, $f_x(x) = e^{-\lambda} \cdot \frac{\lambda^x}{x!}$

CONTINUOUS RV CASE

PROBABILITY DENSITY FUNCTION (PDF)

$$f_x : T_x \rightarrow \mathbb{R}$$

$$P(x \in A) = \int_A f_x(x) dx \quad \xrightarrow{A=[a,b]} \quad P(a \leq x \leq b) = \int_a^b f_x(x) dx$$

$$! P(a) = \int_a^a f_x(x) dx = 0$$

PROPERTIES

$$1) f_x(x) \geq 0 \quad \forall x \in T_x$$

$$2) \int_{x \in T_x} f_x(x) dx = 1$$

3) THE SAME OF DEFINITION

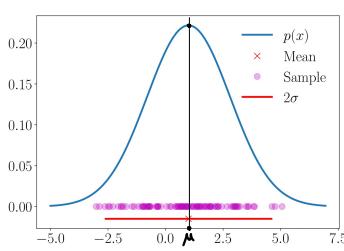
EXAMPLES

1) CONTINUOUS UNIFORM DISTRIBUTION : if $T_x = [a, b]$, then $f_x(x) = \frac{1}{b-a}$

$$-\frac{(x-a)^2}{2\sigma^2}$$

2) NORMAL DISTRIBUTION : $f_x(x) = \frac{1}{\sigma\sqrt{2\pi}} \cdot e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

 $(\mu, \sigma) = (0, 1)$ we have standard normal distribution



10/10

STATISTICS

PROBABILITY → model of some process, where the underlying uncertainty is captured by random variables and we derive what happens

≠

STATISTICS → we observe that something has happened and we figure out the process that explains it

UNIVARIATE DISTRIBUTION

We consider only 1 random variable

BIVARIATE DISTRIBUTION

We consider only 2 random variables

MULTIVARIATE DISTRIBUTION

We consider more random variables

JOINT PROBABILITY

X, Y random variables
 T_x, T_y target spaces
 $T_{xy} = T_x \times T_y$

$f_{xy} : T_{xy} \rightarrow [0, 1]$

$f_{xy}(x_i, y_j) = P(X=x_i, Y=y_j) = \frac{n_{xy}}{N}$

MARGINAL PROBABILITY

$y \setminus x$	x_1	x_2	...	x_k	$P_Y(y)$
y_1	$P(x_1, y_1)$.	.	.	$\sum_i P(x_i, y_1)$
\vdots	\vdots
y_n	.	.	.	$P(x_n, y_n)$	$\sum_i P(x_i, y_n)$
$P_X(x)$	$\sum_j P(x_i, y_j)$...			1

$$P_X(x_i) = \sum_j P(x_i, y_j) \quad \text{MARGINAL PROBABILITY OF } X$$

$$P_Y(y_j) = \sum_i P(x_i, y_j) \quad \text{MARGINAL PROBABILITY OF } Y$$

! THERE IS ALSO FOR DENSITY FUNCTION!

CONDITIONAL PROBABILITY

$$P(Y=y_j | X=x_i) = \frac{n_{ij}}{c_i} \quad \text{column for each } x_i$$

$$P(X=x_i | Y=y_j) = \frac{n_{ij}}{r_j} \quad \text{row for each } y_j$$

SUM RULE

$$p(x) = \begin{cases} \sum_{y \in \mathcal{Y}} p(x, y) & \text{if } y \text{ is discrete} \\ \int_{\mathcal{Y}} p(x, y) dy & \text{if } y \text{ is continuous} \end{cases}$$

PRODUCT RULE

$$p(\mathbf{x}, \mathbf{y}) = p(\mathbf{y} | \mathbf{x})p(\mathbf{x}) = P(y, x)$$

BAYES THEOREM

$$p(\mathbf{x} | \mathbf{y}) = \frac{\underbrace{p(\mathbf{y} | \mathbf{x})}_{\text{posterior}} \underbrace{p(\mathbf{x})}_{\text{prior}}}{\underbrace{p(\mathbf{y})}_{\text{evidence}}} = \frac{p(\mathbf{y} | \mathbf{x})p(\mathbf{x})}{p(\mathbf{y})} = \frac{P(x, y) = P(y|x)P(x)}{P(y, x) = P(x|y)P(y)}$$

↑
← $P(y|x)P(x) = P(x|y)P(y)$

our subjective prior knowledge of x
how x and y are related

We use it to know more about random variables to make inferences on unobservant rv.

QUANTITY OF INTEREST → encapsulates all available information from the prior and the data

↳ instead of considering the posterior in every computation, we consider only some statistic of the posterior (its maximum) → **LOST OF INFORMATION** → if we consider it every time it will lead to a quicker learning

$$p(y) = \int p(y|x)p(x) dx \quad \text{MARGINAL LIKELIHOOD/EVIDENCE}$$

STATISTIC

It's a deterministic function of that random variable.

SUMMARY STATISTIC → useful view of behaviour, summarize and characterize the distribution

EXPECTED VALUE

$$\mathbb{E}_x[g(x)] = \begin{cases} \int_X g(x)p(x) dx & \text{continuous} \\ \sum_{x \in X} g(x)p(x) & \text{discrete} \end{cases}$$

$$\mathbb{E}_X[g(\mathbf{x})] = \begin{bmatrix} \mathbb{E}_{X_1}[g(x_1)] \\ \vdots \\ \mathbb{E}_{X_D}[g(x_D)] \end{bmatrix} \in \mathbb{R}^D \quad \rightarrow \text{MULTIVARIABLE}$$

MEAN

$$\mathbb{E}_{X_d}[x_d] := \begin{cases} \int_{\mathcal{X}} x_d p(x_d) dx_d & \text{continuous} \\ \sum_{x_i \in \mathcal{X}} x_i p(x_d = x_i) & \text{discrete} \end{cases}$$

$$\mathbb{E}_X[\mathbf{x}] = \begin{bmatrix} \mathbb{E}_{X_1}[x_1] \\ \vdots \\ \mathbb{E}_{X_D}[x_D] \end{bmatrix} \in \mathbb{R}^D \quad \rightarrow \text{MULTIVARIABLE}$$

COVARIANCE AND VARIANCE

$$\text{Cov}[\mathbf{x}, \mathbf{y}] = \mathbb{E}[\mathbf{x}\mathbf{y}^\top] - \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{y}]^\top = \text{Cov}[\mathbf{y}, \mathbf{x}]^\top \in \mathbb{R}^{D \times E}$$

If you consider $\mathbf{y} = \mathbf{x}$ $\longrightarrow \text{Cov}(\mathbf{x}, \mathbf{x}) = \mathbb{E}[\mathbf{x}\mathbf{x}^\top] - \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{x}]^\top = \mathbb{V}_{\mathbf{x}} = \sigma^2$

CORRELATION

$$\text{corr}[x, y] = \frac{\text{Cov}[x, y]}{\sqrt{\mathbb{V}[x]\mathbb{V}[y]}} \in [-1, 1]$$

In machine learning, we search correlation to learn better.

EMPIRICAL MEAN

$$\bar{x} := \frac{1}{N} \sum_{n=1}^N x_n$$

EMPIRICAL COVARIANCE

$$\Sigma := \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^\top = \frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})^2$$

SUM AND TRANSFORMATION OF RANDOM VARIABLES

$$\mathbb{E}[x+y] = \mathbb{E}[x] + \mathbb{E}[y]$$

$$\mathbb{E}[x-y] = \mathbb{E}[x] - \mathbb{E}[y]$$

$$\text{V}[x+y] = \text{V}[x] + \text{V}[y] + \text{Cov}[x,y] + \text{Cov}[y,x]$$

$$\text{V}[x-y] = \text{V}[x] + \text{V}[y] - \text{Cov}[x,y] - \text{Cov}[y,x]$$

STATISTICAL INDEPENDENCE

X, Y are statistically independent if $\mathbb{P}(x,y) = \mathbb{P}(x)\mathbb{P}(y)$

PROPERTIES

$$p(y|x) = p(y)$$

$$p(x|y) = p(x)$$

$$\text{V}_{X,Y}[x+y] = \text{V}_X[x] + \text{V}_Y[y]$$

$$\text{Cov}_{X,Y}[x,y] = 0$$

CONDITIONAL INDEPENDENCE

X, Y are conditionally independent $\Leftrightarrow \mathbb{P}(x,y|z) = \mathbb{P}(x|z)\mathbb{P}(y|z)$

INNER PRODUCT OF RANDOM VARIABLES

$\langle X, Y \rangle = \text{Cov}[x,y]$ but X, Y must have $\mu=0$ (mean)

↳ symmetric and positive definite

$$\|x\| = \sqrt{\langle x, x \rangle} = \sqrt{\text{Cov}[x,x]} = \sqrt{\text{V}[x]} = \sigma[x] \quad \text{LENGTH OF RV}$$

$$\cos \theta = \frac{\langle x, y \rangle}{\|x\| \cdot \|y\|} = \frac{\text{Cov}[x,y]}{\sqrt{\text{V}[x] \cdot \text{V}[y]}} \quad \text{ANGLE BETWEEN } X, Y$$

$X \perp Y \Leftrightarrow \langle x, y \rangle = 0 = \text{Cov}[x,y] \quad \text{ORTHOGONALITY}$

GAUSSIAN DISTRIBUTION

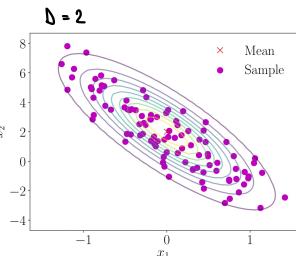
$$p(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

MULTIVARIATE
NOT A PROBABILISTIC DEPENDENCE
BUT JUST A MATHEMATICAL ONE

DIMENSIONS

$$p(x | \mu, \Sigma) = (2\pi)^{-\frac{D}{2}} |\Sigma|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(x - \mu)^\top \Sigma^{-1}(x - \mu)\right)$$

VARIANCE MATRIX



SUMS AND LINEAR TRANSFORMATION

If x, y are independent gaussian variables then $p(x + y) = \mathcal{N}(\mu_x + \mu_y, \Sigma_x + \Sigma_y)$

DATA, MODELS AND LEARNING

GOOD MODEL

Is a model which performs well on unseen data

MODELS

MODELS → FUNCTION (deterministic)
MODELS → PROBABILISTIC MODEL (probability distribution)

MODELS AS FUNCTIONS

We will consider only this class of functions: $f(x) = \theta^\top x + \theta_0$

$$\theta = [\theta_0, \dots, \theta_d]^\top$$

MODELS AS PROBABILITY DISTRIBUTION

We consider data as noisy and we want, by ML, to extract the information from the noise.

We will consider them as multivariate probability distribution.

LEARNING IS FINDING PARAMETERS

The right parameters will make my model a good model,

There are 3 algorithmic phases:

- 1) PREDICTION OR INFERENCE
- 2) TRAINING OR PARAMETER ESTIMATION
- 3) HYPERPARAMETER TUNING OR MODEL SELECTION

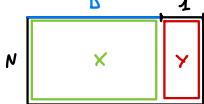
EMPIRICAL RISK MINIMIZATION

It's the learning phase for non-probabilistic model

HYPOTHESIS CLASS OF FUNCTION

$x_n \in \mathbb{R}^D$ values and $y_n \in \mathbb{R}^M$ labels, by supervised learning setting we have $(x_1, y_1), \dots, (x_N, y_N)$

So we have an $N \times (D+1)$ matrix



And we hope to find a good parameter θ^* $\rightarrow f(x_n, \theta^*) \approx y_n$

\hat{y}_n

"

LOSS FUNCTION FOR TRAINING

$$\ell(y_n, \hat{y}_n)$$

We assume $(x_1, y_1), \dots, (x_N, y_N)$ as independent and identically distributed

$$X = [x_1, \dots, x_N]^T \in \mathbb{R}^{N \times D}$$
 and $y = [y_1, \dots, y_N]^T$

$$R_{\text{emp}}(f, X, y) = \frac{1}{N} \sum_{n=1}^N \ell(y_n, \hat{y}_n)$$

IN LAB WE CONSIDERED ALL THIS AS LOSS FUNCTION

$$\text{ES: } \ell(y_n, \hat{y}_n) = (y_n - \hat{y}_n)^2 \quad \text{LEAST-SQUARES LOSS}$$

$$\text{we want } \min \{\ell(y_n, \hat{y}_n)\} = \theta^*$$

$$\min_{\theta \in \mathbb{R}^D} \frac{1}{N} \sum_{n=1}^N (y_n - f(x_n, \theta))^2 \rightarrow \min_{\theta \in \mathbb{R}^D} \frac{1}{N} \sum_{n=1}^N (y_n - \theta^\top x_n)^2 \rightarrow \min_{\theta \in \mathbb{R}^D} \frac{1}{N} \|y - X\theta\|^2 = \theta^*$$

$n = N+1, \dots, M$

We want also to minimize the EXPECTED RISK $R_{\text{true}}(f) = \mathbb{E}_{x,y}[\ell(y, f(x))]$

REGULARIZATION TO REDUCE OVERFITTING

Overfitting happens when the predictor fits too closely the training data and doesn't generalize well to new data. $\hookrightarrow R_{\text{EMP}}(f, X_{\text{TRAIN}}, Y_{\text{TRAIN}})$ underestimates $R_{\text{TRUE}}(f)$

We find a compromise by using regularization by adding a REGULARIZATION PARAMETER and REGULATOR

$$\text{ES: } \ell(y_n, \hat{y}_n) = (y_n - \hat{y}_n)^2 \quad \text{LEAST-SQUARES LOSS}$$

$$\theta^* = \min_{\theta} \frac{1}{N} \|y - X\theta\|^2 + \lambda \|\theta\|^2$$

PARAM \uparrow REGULATOR
PENALTY TERM

PARAMETER ESTIMATION

We will use probability distribution to model our uncertainty

MAXIMUM LIKELIHOOD ESTIMATION (MLE)

- LIKELIHOOD \rightarrow function related to probability distribution $\mathcal{L}_{\underline{x}}(\theta) = -\log(p(x | \theta))$

If we assume $(x_1, y_1), \dots, (x_n, y_n)$ independent and identically distributed so we have that
 $p(\mathcal{Y} | \mathcal{X}, \theta) = \prod_{n=1}^N p(y_n | x_n, \theta) \longrightarrow \mathcal{L}(\theta) = -\sum_{n=1}^N \log(p(y_n | x_n, \theta))$ TO FIND GOOD WE MINIMIZE $L(\theta)$

ES: GAUSSIAN DISTRIBUTION

$y_n = x_n^\top \theta$ for prediction and $\varepsilon_n \sim N(0, \sigma^2)$ OBSERVATION UNCERTAINTY

So we have that $p(y_n | x_n, \theta) = N(y_n | x_n^\top \theta, \sigma^2) \longrightarrow$

$$\min_{\theta} \left\{ \|y - X\theta\|_2^2 \right\} \leftarrow \min_{\theta} \left\{ \sum_{n=1}^N (y_n - x_n^\top \theta)^2 \right\} \xleftarrow{\text{MINIMIZE } L(\theta)}$$

MAY SUFFER OF OVERFITTING

$$\begin{aligned} \mathcal{L}(\theta) &= -\sum_{n=1}^N \log(p(y_n | x_n, \theta)) = -\sum_{n=1}^N \log(N(y_n | x_n^\top \theta, \sigma^2)) \\ &= -\sum_{n=1}^N \log\left(\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_n - x_n^\top \theta)^2}{2\sigma^2}\right)\right) \\ &= -\sum_{n=1}^N \log\left(\exp\left(-\frac{(y_n - x_n^\top \theta)^2}{2\sigma^2}\right)\right) - \sum_{n=1}^N \log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) \\ &= \frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - x_n^\top \theta)^2 - \sum_{n=1}^N \log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right). \end{aligned}$$

LEAST SQUARE CONSTANT

WE CAN JUST MINIMIZE LEAST SQUARE

MAXIMUM A POSTERIORI ESTIMATION (MAP)

If we have a prior knowledge x that makes us having more knowledge on θ .

$$p(\theta | x) = \frac{p(x | \theta)p(\theta)}{p(x)} \longrightarrow p(\theta, x) \propto p(x, \theta) \downarrow p(\theta)$$

LINEAR REGRESSION

We want to find a function f that maps $x \in \mathbb{R}^d$ in $f(x) \in \mathbb{R}$ with an error ϵ

$$y_n = f(x_n) + \epsilon$$

In order to do that:

1) CHOICE OF MODEL AND PARAMETRIZATION

Given a dataset, which function classes can best model the data?

2) FINDING GOOD PARAMETERS

Given the model, which are the parameters that minimize the loss?

3) OVERFITTING AND MODEL SELECTION

Given too good parameters, how much error do I need to introduce to make the model more general

$$y_n = f(x_n) + \epsilon \quad \text{and} \quad P(y|x) = N(y|f(x), 0)$$

if $f = x^T \theta$ it's the previous case

NONLINEAR TRANSFORMATION OF x → BECAUSE BEFORE IT WAS UNDERFITTING

We suppose that $f(x_n) = \Xi^T(x) \cdot \theta$

$$\Xi^T(x) \cdot \theta = \sum_{n=1}^N \Xi(x_n) \theta_n \quad \Xi : \mathbb{R}^d \rightarrow \mathbb{R}^k \quad \text{and} \quad \Xi_i : \mathbb{R}^d \rightarrow \mathbb{R}$$

DSK
USUALLY

ES: POLYNOMIAL REGRESSION

$$\text{We have } \Xi_i(x) = x^i \quad \text{and so} \quad f(x) = \sum_{i=0}^{k-1} \Xi_i(x) \theta_i = \sum_{i=0}^{k-1} x^i \theta_i$$

THE PREDICTIONS ARE NOT LINES ANYMORE BUT THEY ARE FUNCTIONS OF SETS OF POSSIBLE POLYNOMIAL

To solve it (like $f(x) = x^T \theta$) we have $\ell(\theta) = \frac{1}{2\sigma^2} \|Y - \Xi\theta\|^2$ to be minimized so we use gradient $\nabla(\|Y - \Xi\theta\|^2) = \emptyset$

$$\nabla(\|Y - \Xi\theta\|^2) = \Xi^T \Xi \theta - \Xi^T Y = \emptyset \quad \longrightarrow \quad \theta^* = \Xi^T Y \cdot (\Xi^T \Xi)^{-1}$$

17/10

POLYNOMIAL REGRESSION

1) DATA SET $(x_i, y_i) \rightarrow y_i = f(x_i) + \epsilon_i \quad \epsilon_i \sim N(0, \sigma^2)$

2) $\hat{\theta}^* = \min_{\theta} \{ \|Y - \Phi\theta\|_2^2\}$ by SGD

3) $f(x) = \sum_{i=0}^{k-1} \theta_i \cdot x^i$ polynomial

4) MAP $\rightarrow \hat{\theta}_{MAP}^* = \min_{\theta} \{ \|Y - \Phi\theta\|_2^2 + \lambda \|\theta\|_2^2\}$ by SGD

5) $f(x) = \sum_{i=0}^{k-1} \hat{\theta}_{MAP}^* \cdot x^i$

RMSE (Root Mean Square Error)

$$RMSE = \sqrt{\frac{1}{N} \|y - \Phi\theta\|^2} = \sqrt{\frac{1}{N} \sum_{n=1}^N (y_n - \phi^\top(x_n)\theta)^2}$$

COMPARE ERRORS ON DATASET OF DIFFERENT DIMENSION

HAS SAME SCALE AND UNITS AS y_n

NOT OVERFITTING ANYMORE

INCREASING THE EXPONENTIAL
OF THE POLYNOMIAL WILL
INTRODUCE MORE ERROR

